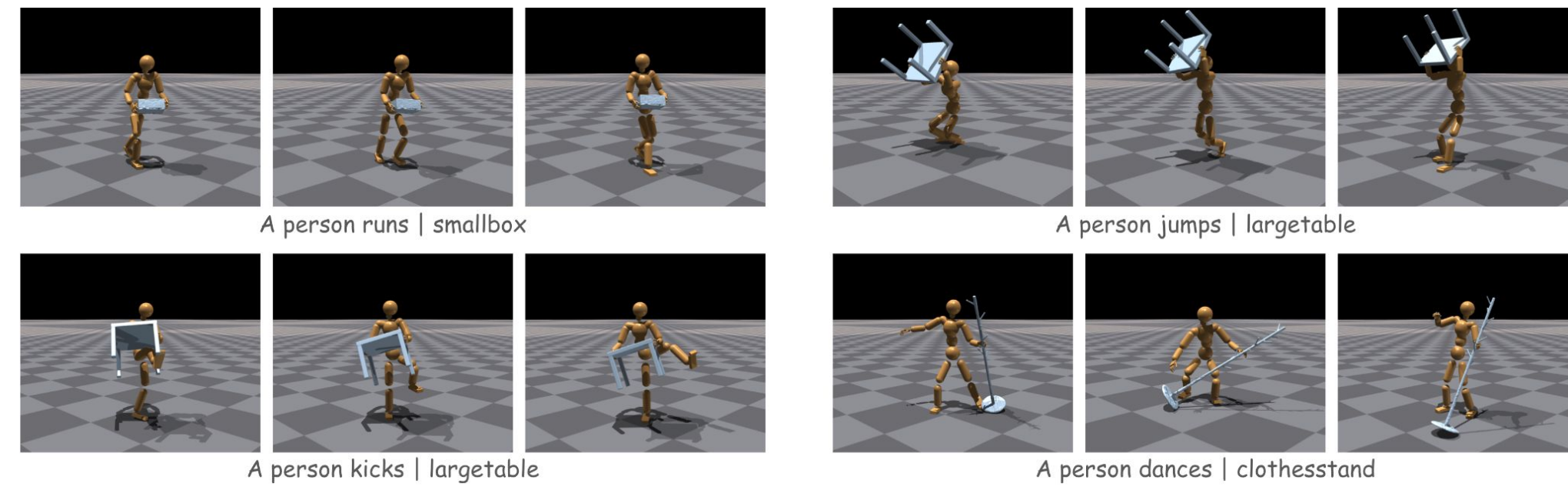




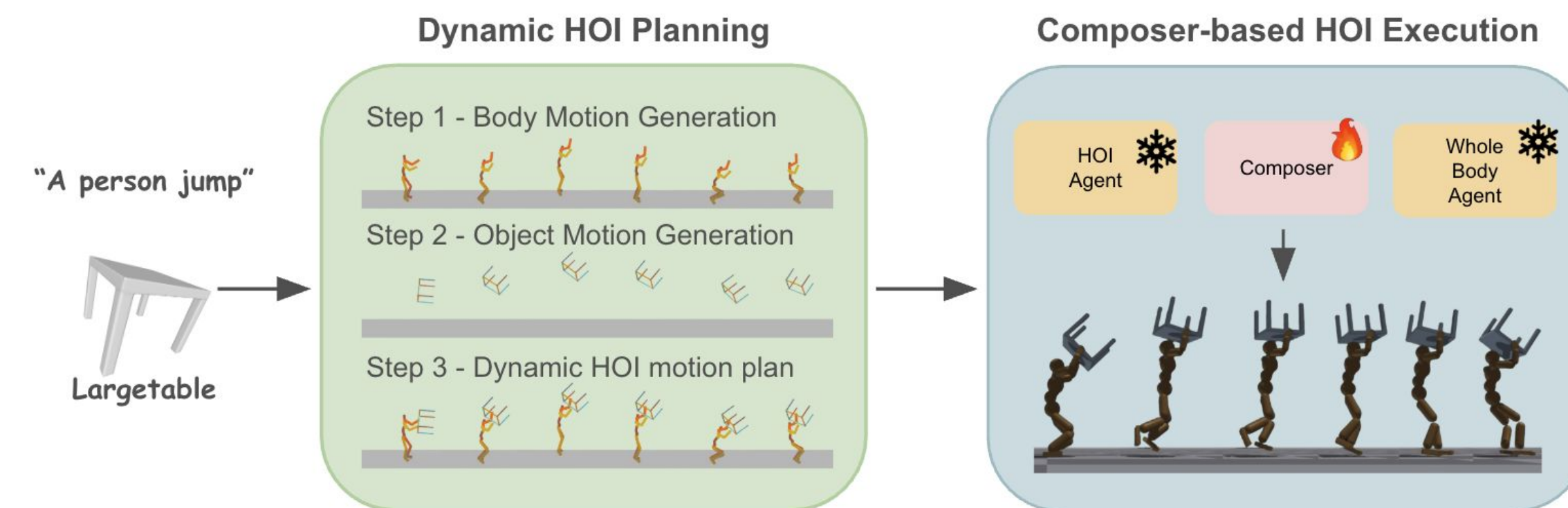
## Introduction



run, jump, kick, dance while interacting with an object

- Plan and execute physically plausible, **dynamic** and **contact-rich Human-Object Interaction (HOI)**
- Existing datasets and agents handle either HOI or dynamic motions, but **not both simultaneously**

## Overview



- Dynamic HOI Planning:** Plan dynamic motions with consistent hand-object contact
- Composer-based HOI Execution:** Efficiently learn dynamic HOI control by blending pre-trained experts

## Quantitative Results

Method	HOI	Physics	HOI QUALITY		PHYSICAL PLAUSIBILITY				MOTION QUALITY	
			$C_{\%} \uparrow$	$C_{cons} \uparrow$	$Pene_{obj} \downarrow$	Skate $\downarrow$	Float $\downarrow$	Jitter <sub>pos</sub> $\downarrow$	R-Prec $\uparrow$	Diversity $\uparrow$
MDM			-	-	-	<b>0.133</b>	29.3	$9.43 \times 10^4$	<b>0.374</b>	<b>7.61</b>
HOI-Diff	✓		0.285	18.2	8.18	0.633	<b>17.8</b>	$1.56 \times 10^4$	0.257	4.66
DAViD	✓		0.848	10.9	4.842	0.261	20.4	$6.69 \times 10^4$	0.310	6.70
Ours <sub>p</sub>	✓		<b>1.000</b>	<b>0.906</b>	4.196	<b>0.217</b>	30.1	$5.93 \times 10^4$	0.332	5.56
Ours <sub>p+E</sub>	✓	✓	<b>0.999</b>	2.95	<b>0.009</b>	0.786	<b>9.95</b>	$4.53 \times 10^4$	0.316	3.97

**Generation Quality**  
(Ours<sub>p</sub>: planning only / Ours<sub>p+E</sub>: full planning-execution pipeline)

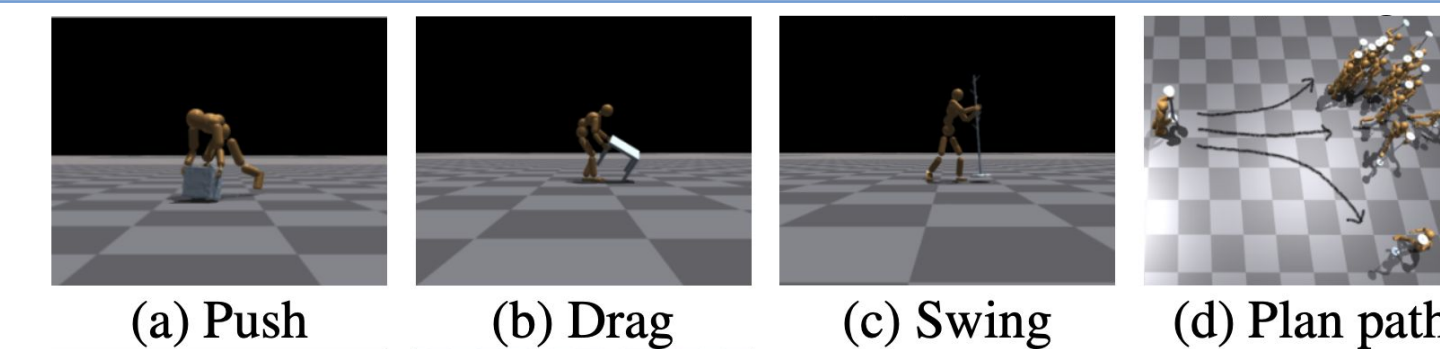
- Ours<sub>p</sub> achieves the **best HOI Quality** –  $C_{\%}$  and  $C_{cons}$
- Ours<sub>p+E</sub> achieves the **best physical plausibility** –  $Pene_{obj}$  and Float, 2nd Jitter<sub>pose</sub>

Method	SR $\uparrow$	T	D $\uparrow$	$E_{HOI} \downarrow$	Jitter <sub>DoF</sub> $\downarrow$	Blending			
						SR $\uparrow$	D $\uparrow$	$E_{HOI} \downarrow$	
PPO	0.227	75	0.859	59.098	$0.424 \times 10^3$	Heuristic <sub>Hand</sub>	0.365	1.069	13.04
PHC [26]	0.397	-	1.166	12.965	$1.350 \times 10^3$	Heuristic <sub>Arm</sub>	0.472	1.603	17.18
PHC <sub>R</sub>	0.313	15	1.834	26.024	$1.762 \times 10^3$	Hard MoE	0.383	2.748	20.48
InterMimic [56]	0.376	-	2.407	21.379	$0.436 \times 10^3$	Hard MoE (Joint)	0.383	2.535	26.45
InterMimic <sub>R</sub>	0.310	15	2.324	34.675	<b><math>0.336 \times 10^3</math></b>	Ours <sub>MLP</sub>	0.571	3.034	18.03
InterMimic <sub>FT</sub>	<b>0.526</b>	75	<b>7.37</b>	<b>8.635</b>	$0.370 \times 10^3$	Ours <sub>MLP+PCA</sub>	<b>0.591</b>	<b>4.607</b>	<b>11.67</b>
Ours	<b>0.591</b>	23	<b>4.607</b>	<b>11.667</b>	$0.359 \times 10^3$				

**Execution Performance**  
(Ours<sub>MLP</sub>: except PCA subspace exploration)

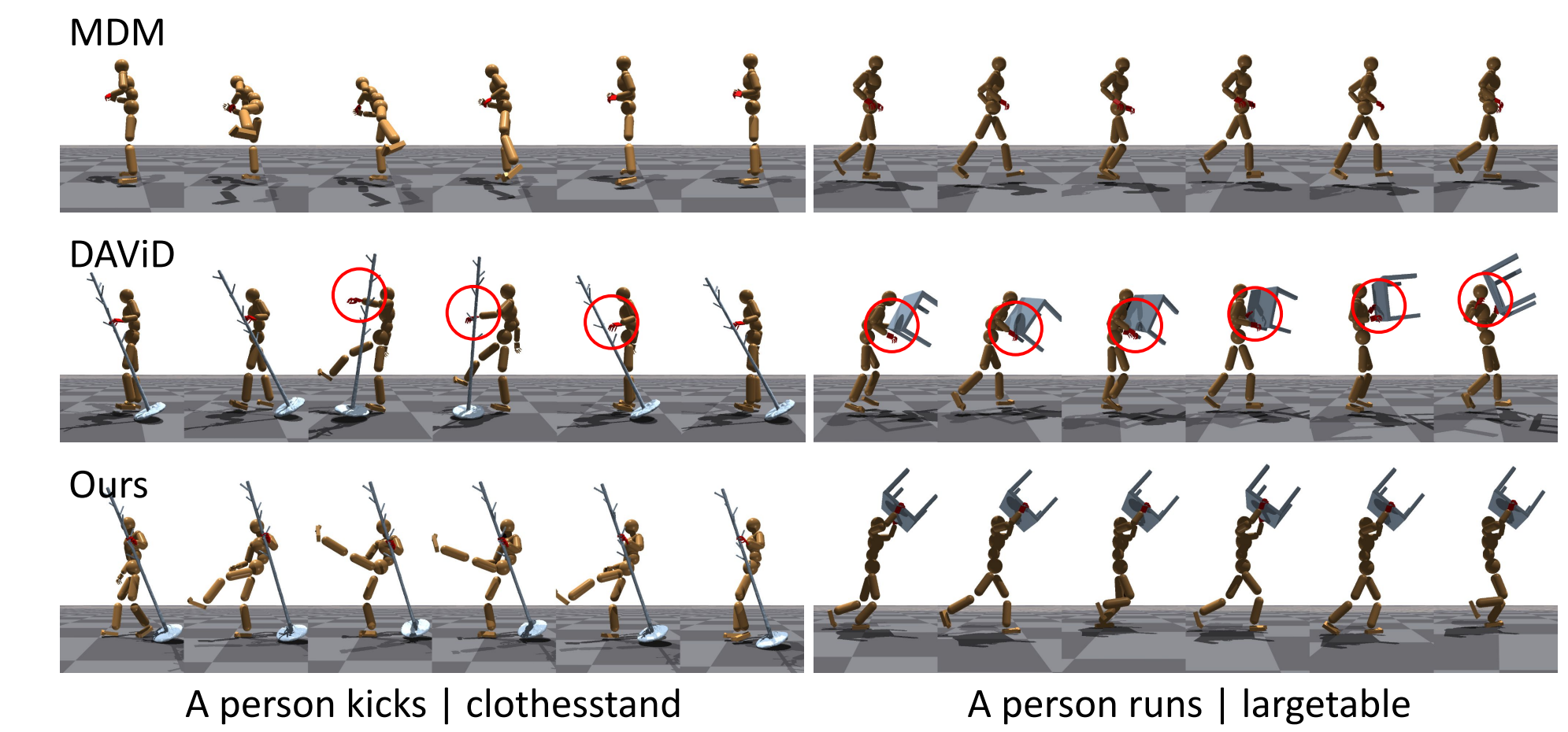
- Ours achieves the **highest success rate**
- Converges **3x faster** than fine-tuning InterMimic
- Ours blending method achieved best SR, D, and  $E_{HOI}$

## Scalability



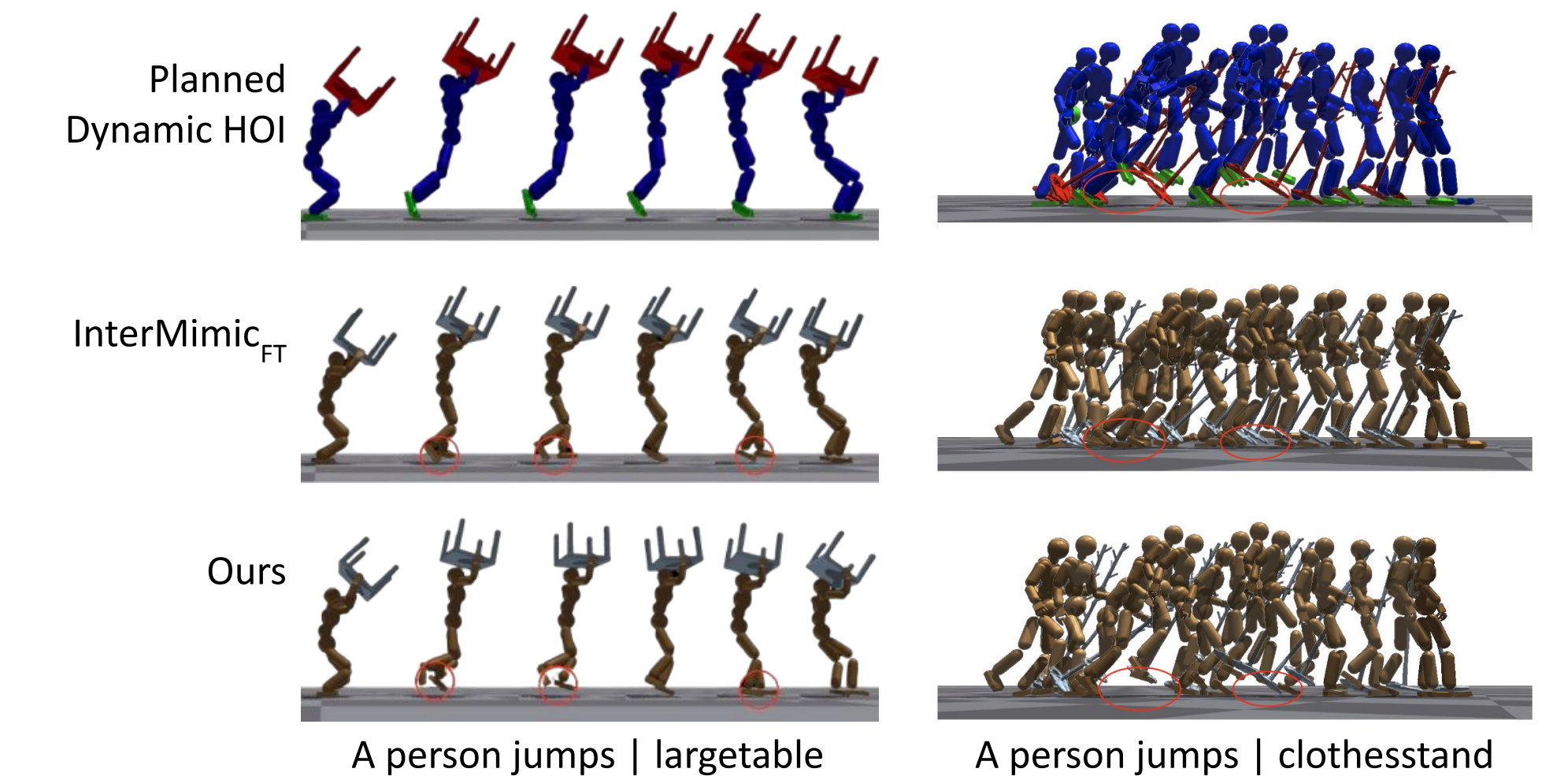
- Diverse object geometries (FullBodyManip)
- Diverse motion styles (push, drag, swing)
- Plug-and-play application – e.g., MDM → GMD enables path planning

## Qualitative Results



Comparison of HOI Planning

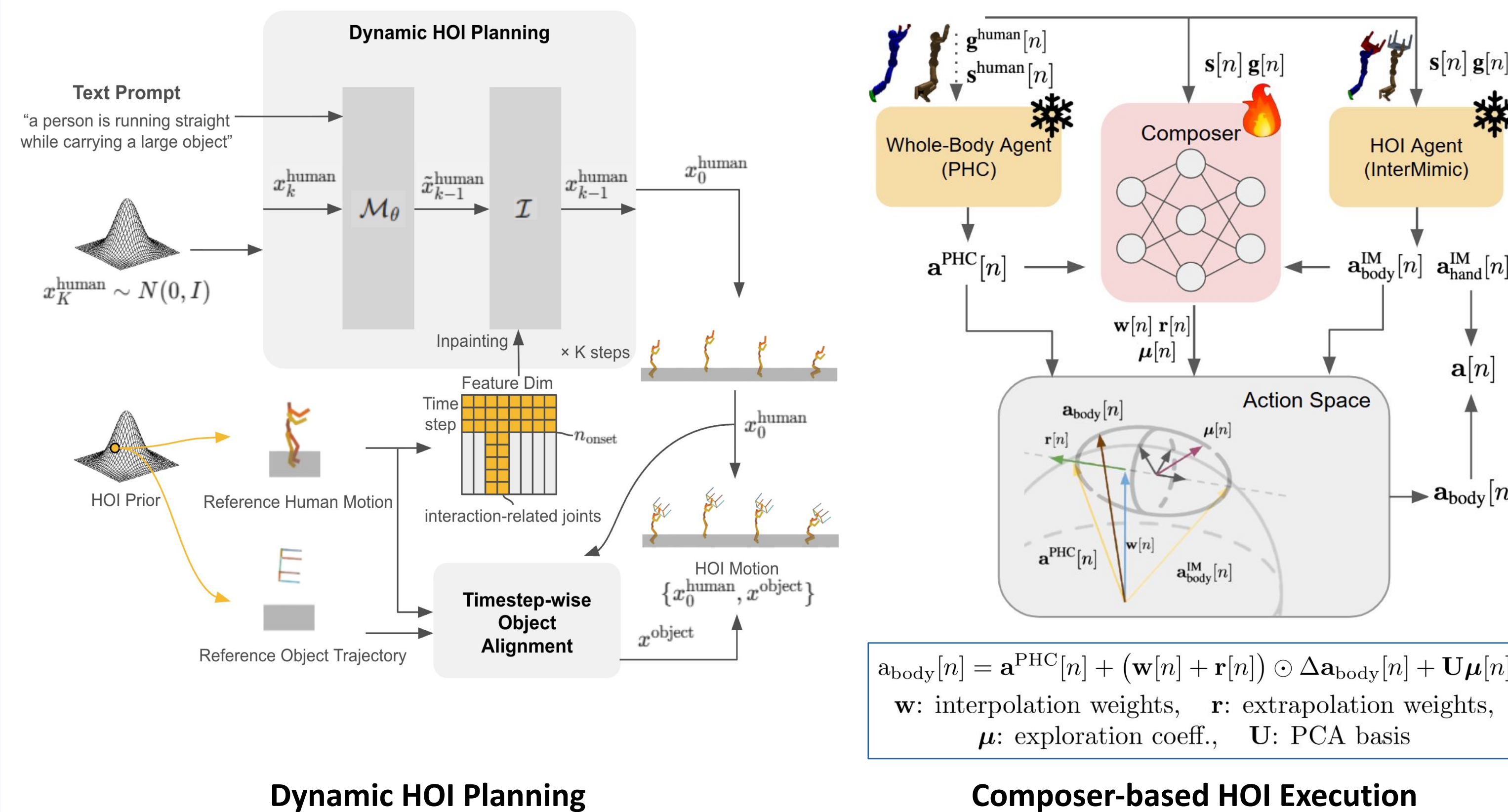
- Ours accurately follows the text-aligned **dynamic motion** while maintaining **consistent hand-object contact**



Comparison of HOI Execution

- Ours successfully executes the planned dynamic HOI while maintaining **motion style** and **object interaction**
- InterMimic<sub>FT</sub> collapses to a local minimum – shuffling instead of two-foot takeoff

## Methodology



### Dynamic HOI Planning

- Inpaint HOI prior (FullBodyManip) into MDM during diffusion sampling
- Fix interaction-related joints after contact onset; text-aligned motion inpainting for the others
- Compute object trajectories from hand poses

### Composer-based HOI Execution

- Composer predicts per-DoF inter-/extrapolation weights
- PCA-based low-dim subspaces for further exploration

$$\mathbf{a}_{body}[n] = \mathbf{a}^{PHC}[n] + (\mathbf{w}[n] + \mathbf{r}[n]) \odot \Delta \mathbf{a}_{body}[n] + \mathbf{U} \boldsymbol{\mu}[n]$$

$\mathbf{w}$ : interpolation weights,  $\mathbf{r}$ : extrapolation weights,  $\boldsymbol{\mu}$ : exploration coeff.,  $\mathbf{U}$ : PCA basis

## Contact

Sanghyeok Nam  
KAIST, Republic of Korea  
Email: sang990701@kaist.ac.kr  
Website: <https://www.linkedin.com/in/shnam99>

## References

- N. Mahmood et al., "AMASS: Archive of Motion Capture as Surface Shapes," ICCV, 2019.
- J. Li et al., "Object Motion Guided Human Motion Synthesis," ACM TOG, 2023.
- G. Tevet et al., "Human Motion Diffusion Model," ICLR, 2023.
- X. Peng et al., "HOI-Diff: Text-Driven Synthesis of 3D Human-Object Interactions Using Diffusion Models," CVPR Workshop on HuMoGen, 2025.
- H. Kim et al., "DAViD: Modeling Dynamic Affordance of 3D Objects Using Pretrained Video Diffusion Models," ICCV, 2025.
- Z. Luo et al., "Perpetual Humanoid Control for Real-time Simulated Avatars," ICCV, 2023.
- S. Xu et al., "InterMimic: Towards Universal Whole-Body Control for Physics-Based Human-Object Interactions," CVPR, 2025.
- K. Karunatanakul et al., "Guided Motion Diffusion for Controllable Human Motion Synthesis," ICCV, 2023.